



## Error estimates using the cell discretization method for steady-state convection–diffusion equations

Howard Swann\*

*San José State University, San José, CA 95192-0103, USA*

Received 25 July 1996; received in revised form 15 February 1997

### Abstract

The cell discretization algorithm is applied to generate approximate solutions for some second-order non-self-adjoint elliptic equations. General convergence for homogeneous problems is shown by obtaining suitable error estimates. The method is applied using polynomial bases; this provides a nonconforming extension of the finite element method that can also produce the continuous approximations of an  $h$ - $p$  finite element method. Numerical tests on convection–diffusion problems are made that confirm the theoretical estimates, and methods for dealing with boundary layer problems are illustrated.

**Keywords:** Elliptic second-order partial differential equations; Non-self-adjoint; Nonconforming finite element methods; Primal hybrid method; Cell discretization; Convection–diffusion equation

**AMS classification:** primary 65N30; 65M15 secondary 65N12

### 1. Introduction

A second-order linear partial differential operator is of the form

$$Au \equiv \sum_{i,j=1}^K b_{ij} D_i D_j u + \sum_{i=1}^K b_i D_i u + a_0 u,$$

where  $D_i$  denotes partial differentiation with respect to  $x_i$ . By absorbing certain derivatives of the functions  $b_{ij}(x)$  into the first-order terms  $\sum_{i=1}^K b_i D_i u$ , we can convert this operator to the standard form

$$Au \equiv - \sum_{i,j=1}^K D_i (a_{ij} D_j u) + \sum_{i=1}^K a_i D_i u + a_0 u$$

\* Corresponding author. E-mail: swann@sjsumcs.sjsu.edu.

by defining appropriate  $a_{ij}$  that will satisfy  $a_{ij}(x) = a_{ji}(x)$ . The operator is *elliptic* if there exists  $c > 0$  such that

$$\sum_{i,j}^K a_{ij}(x) z_i z_j \geq c \sum_{i=1}^K z_i^2$$

for  $x$  in domain  $\Omega$  for any  $z_i \in \mathbb{R}$ . The operator is said to be *self-adjoint* if the first-order terms  $\sum_{i=1}^K b_i D_i u$  are absent. The standard example for a self-adjoint elliptic problem is the Poisson equation

$$-\Delta u = f.$$

Self-adjoint elliptic problems can be converted to a variational setting, where a certain functional is to be minimized. When applied to an approximation space consisting of finite linear combinations of basis functions, the functional induces a quadratic form and standard procedures for finding a minimum of a quadratic form produce a system of linear equations whose solutions provide the coefficients for the approximation. However, such variational methods do not readily apply to non-self-adjoint equations; thus some other method is more appropriate.

Our results hold for general non-self-adjoint second-order elliptic equations; the standard example is given by the steady-state convection–diffusion equations

$$\nabla \cdot (-\sigma \nabla u + u \mathbf{a}) = 0,$$

which describe the concentration  $u(x)$  of a substance at point  $x$ , where the fluid that carries the substance is subject to a steady-state flow with velocity  $\mathbf{a}(x)$  and  $\sigma$  is the diffusion coefficient.

We use a nonconforming extension of the finite element method called the *cell discretization algorithm* by Greenstadt [7–10]. It generalizes the *primal hybrid method* of Raviart and Thomas [14]. A domain is partitioned into *cells*; approximate solutions are given by linear combinations of functions that are part of any Schauder basis on each cell. We enforce a form of weak continuity on approximations over the entire domain by requiring that the difference of the traces of approximations on the common boundaries of adjacent cells be orthogonal to selected functions that are part of a basis defined on the interfaces between cells. These requirements, called *moment collocations*, give a set of linear constraints on the coefficients of the basis functions on each cell. For self-adjoint problems, the usual variational formulation allows us to use Lagrange multipliers to enforce these collocation constraints [4, 17]. This approach is not possible here, but we can use methods for solving parabolic equations [18] to construct a basis that carries with it the weak continuity constraints and then follows the standard Galerkin approximation procedure.

We describe the algorithm in Section 2 and, applying the results in [4, 17], we obtain error estimates that show convergence of approximations in  $H^1$  on each cell to solutions of a homogeneous boundary value problems as the number of moment collocations enforced becomes large and the number of basis functions utilized becomes suitably larger.

In Section 3 we give methods for nonhomogeneous problems and convert our general error estimates to a polynomial implementation of the algorithm for domains partitioned into triangles and parallelograms. We express our error estimates in terms of an “ $h$ – $p$ ” finite element context [2, 3, 12]. We have created appropriate software that can generate the continuous approximations of a finite element method for comparison with our nonconforming approximations.

Our methods provide convergent approximations when applied to steady-state convection–diffusion equations when the underlying flow is constant or incompressible. Of particular interest are convection–dominated convection–diffusion equations. Two examples are given. A good approximation is obtained in the first example from [15]. The second example similar to those in [13] shows how a good choice of partition can assist in dealing with a boundary layer.

## 2. Convergence results

The setting for the method is described in full in [17, 18]. We summarize here.

We assume that bounded domain  $\Omega$  in  $\mathbb{R}^K$  (with boundary  $\Gamma$ ) is partitioned into  $N$  subdomains  $\Omega_i$  with Lipschitz continuous boundaries that are piecewise  $C^1$ ; such subdomains are called *cells*. The exterior is  $\Omega_0 \equiv \mathbb{R}^K \setminus \overline{\Omega}$ .

Let  $H^1(\Omega_i)$  denote the usual Hilbert space on each cell  $\Omega_i$  with inner product  $(\cdot, \cdot)_{1,i}$ . The  $H^1(\Omega_i)$  norm is denoted  $\|\cdot\|_{1,i}$ .  $(\cdot, \cdot)_i$  represents the  $L_2(\Omega_i)$  inner product, with the norm expressed as  $\|\cdot\|_{0,i}$ .  $(\cdot, \cdot)$  denotes the  $L_2(\Omega)$  inner product.

These spaces are assembled to form Hilbert space

$$H \equiv \{u \in L_2(\Omega): u|_{\Omega_i} \in H^1(\Omega_i); i = 1, \dots, N\}$$

with inner product

$$(u, v)_H \equiv \sum_{i=1}^N (u, v)_{1,i}$$

and norm represented by  $\|\cdot\|_H$ .

Let  $\Gamma_{ij}$  represent  $\overline{\Omega_i} \cap \overline{\Omega_j}$ .  $\Gamma_{i0}$  is a boundary segment between  $\Omega_i$  and  $\Omega_0$ . The inner product for  $L_2(\Gamma_{ij})$  is denoted by  $\langle \cdot, \cdot \rangle_{ij}$ , with norm represented as  $\|\cdot\|_{ij}$ .

We denote by  $\gamma_{ij}$  the trace operator restricting  $u|_{\Omega_i}$  to its values on  $\Gamma_{ij}$ . There are constants  $C_{ij}$  that depend on the geometry of  $\Omega_i$  such that for any  $u \in H$ ,  $\|\gamma_{ij}(u)\|_{ij} \leq C_{ij} \|u\|_{1,i}$ .

For each  $\Gamma_{ij}$ , choose  $\{\omega_k^{ij}\}_{k=1}^\infty$  to be functions in  $H^{1/2}(\Gamma_{ij})$  that are a Schauder basis for  $L_2(\Gamma_{ij})$ . Thus, for any  $h \in L_2(\Gamma_{ij})$ , there are coefficients  $d_k$  such that  $h = \sum_{k=1}^\infty d_k \omega_k^{ij}$ . Let  $\mathcal{T}_n^{ij}(h) \equiv \sum_{k=n+1}^\infty d_k \omega_k^{ij}$ . For any  $h$  and  $\varepsilon > 0$ , there is some  $N(h, \varepsilon)$  such that  $n > N(h, \varepsilon) \Rightarrow \|\mathcal{T}_n^{ij}(h)\|_{ij} < \varepsilon$ .

For any  $u \in H$ , we define the  $k$ th *moment* of  $u|_{\Omega_i}$  on  $\Gamma_{ij}$  to be

$$M_k^{ij}(u) \equiv \langle \gamma_{ij}(u), \omega_k^{ij} \rangle_{ij}.$$

Let  $N_I$  be the number of interfaces  $\Gamma_{ij}$ . Form multi-index  $[n]$ , an  $N_I$ -vector of nonnegative integers  $(\dots, n_{ij}, \dots)$ , with integer  $n_{ij}$  associated with interface  $\Gamma_{ij}$ .

Let  $G[n] \equiv \{u \in H: \text{for any } ij, j \neq 0, ij = 1, \dots, N_I \text{ and for any } k \leq n_{ij}, \text{ we have } M_k^{ij}(u) = M_k^{ji}(u)\}$ ; this is the set of functions  $u$  in  $H$  such that the difference of the traces from either side of any internal interface  $\Gamma_{ij}$ ,  $\gamma_{ij}(u) - \gamma_{ji}(u)$ , is  $L_2(\Gamma_{ij})$ -orthogonal to  $\omega_k^{ij}$ ,  $k = 1, \dots, n_{ij}$ . We call such weak continuity across interfaces *moment collocation*.

Let  $G_0[n] \equiv \{u \in G[n]: \text{for any } i, \text{ for any } k \leq n_{i0}, M_k^{i0}(u) = 0\}$ .  $G_0[n]$  is the set of functions in  $G[n]$  that are weakly zero on the external interfaces  $\Gamma_{i0}$  making up  $\Gamma$ .

We define a partial order for such multi-indices; we say  $[n'] \geq [n] \Leftrightarrow$  for any  $ij$ ,  $n'_{ij} \geq n_{ij}$ . If  $[n^k]$  is a sequence of multi-indices,  $k = 1, 2, \dots$ , we say that  $[n^k] \rightarrow [\infty]$  if  $[n^k] \leq [n^{k+1}]$  and  $\inf_{ij} \{n^k_{ij}\} \rightarrow \infty$  as  $k \rightarrow \infty$ .

For the  $i$ th cell  $\Omega_i$ , choose any Schauder basis  $\{B_k^i\}$  for  $H^1(\Omega_i)$ . Thus, for any  $v$  in  $H^1(\Omega_i)$ , there are  $b_k^i$  such that  $\sum_{k=1}^{\infty} b_k^i B_k^i = v$ ; let  $v_{\cdot, m} = \sum_{k=1}^m b_k^i B_k^i$ . Let  $\mathcal{Q}_m^i(v)$  denote the orthogonal projection (in the  $H^1(\Omega_i)$  inner product) of  $v$  onto the  $H^1(\Omega_i)$ -orthogonal complement of the span of  $\{B_1^i, B_2^i, \dots, B_m^i\}$ . Thus

$$\mathcal{Q}_m^i(v_{\cdot, m}) = 0; \quad \mathcal{Q}_m^i(v) = \mathcal{Q}_m^i(v - v_{\cdot, m});$$

$$\|\mathcal{Q}_m^i(v)\|_{1,i} \leq \|v - v_{\cdot, m}\|_{1,i} = \left\| \sum_{k=m+1}^{\infty} b_k^i B_k^i \right\|_{1,i} \quad \text{and} \quad \lim_{m \rightarrow \infty} \|\mathcal{Q}_m^i(v)\|_{1,i} = 0.$$

Let  $[m]$  be an  $N$ -dimensional multi-index indicating the number of basis functions used in the approximation on each of the  $N$  cells; we adopt the same notational conventions as those used for multi-index  $[n]$ .

Define  $H[m]$  to be the subspace of  $H$  such that for any  $v \in H[m]$ ,  $v|_{\Omega_i}$  is in the span of  $\{B_1^i, B_2^i, \dots, B_{m_i}^i\}$ .

Given  $[m]$ , and any function  $v$  in  $H$ ,  $\mathcal{Q}_{[m]}(v)$  is defined to be the function in  $H$  such that  $\mathcal{Q}_{[m]}(v)|_{\Omega_i} = \mathcal{Q}_{m_i}^i(v|_{\Omega_i})$ . Thus,  $\mathcal{Q}_{[m]}(\cdot)$  is the projection of  $H$  onto  $H[m]^\perp$ ;  $\lim_{[m] \rightarrow [\infty]} \|\mathcal{Q}_{[m]}(v)\|_H = 0$ .

Let  $G_0[n][m] \equiv G_0[n] \cap H[m]$ . The moment collocation requirements are met by requiring that certain linear equations hold among the  $b_k^i$ , e.g., for  $u \in G_0[n][m]$ , we require that, on internal interfaces  $\Gamma_{ij}$ ,

$$\langle \gamma_{ij}(u), \omega_p^{ij} \rangle_{ij} - \langle \gamma_{ji}(u), \omega_p^{ij} \rangle_{ij} = 0, \quad p = 1, \dots, n_{ij},$$

which gives the requirement

$$\sum_{k=1}^{m_i} b_k^i \langle \gamma_{ij}(B_k^i), \omega_p^{ij} \rangle_{ij} - \sum_{k=1}^{m_j} b_k^j \langle \gamma_{ji}(B_k^j), \omega_p^{ij} \rangle_{ij} = 0 \quad (1)$$

and, for the external boundary segments  $\Gamma_{i0}$ ,

$$\langle \gamma_{i0}(u), \omega_p^{i0} \rangle_{i0} = 0, \quad p = 1, \dots, n_{i0},$$

which becomes

$$\sum_{k=1}^{m_i} b_k^i \langle \gamma_{ij}(B_k^i), \omega_p^{i0} \rangle_{i0} = 0. \quad (2)$$

The following lemma connects estimates in terms of these spaces.

**Lemma 1.1.** *If  $\mathcal{P}_m^n$  is the orthogonal projection operator of  $G_0[n]$  onto  $G_0[n][m]$ , then there is a constant  $K = K([n])$  depending on the number ‘ $[n]$ ’ of moment collocations, the cell decomposition of domain  $\Omega$  and the choice of basis functions and collocation functions such that*

$$\|u - \mathcal{P}_m^n u\|_H \leq K([n]) \|\mathcal{Q}_{[m]} u\|_H,$$

where  $\mathcal{Q}_{[m]}$  is the projection of  $H$  onto the  $H$ -orthogonal complement of  $H[m]$ .

The proof of this lemma is found in [17]; see [4] for a detailed discussion of – and estimates for –  $K([n])$ . We return to the consideration of  $K([n])$  in the discussion of a polynomial implementation of the algorithm given in Section 3.

We construct a new basis for  $G_0[n][m]$  that carries with it the moment collocation constraints.

The coefficients  $\{b_k^i\}$  for the representation on each of the  $N$  cells can be concatenated to form vector  $\mathbf{b} \equiv (b_1^1, b_2^1, \dots, b_1^2, b_2^2, \dots, b_1^k, b_2^k, \dots)$ . The linear moment collocation requirements (1) and (2) are expressed as  $\mathbf{M}\mathbf{b}^T = \mathbf{0}$ , for a suitable rectangular matrix  $\mathbf{M}$ ; it is an  $n' \times m$  matrix, where  $n' = \sum n_{ij}$  and  $m = \sum_{k=1}^N m_k$ ;  $m > n'$ . It is shown in [17] that, for any  $[n]$ , the rows of  $\mathbf{M}$  are independent if  $[m]$  is sufficiently large. Dorr [6] discusses this phenomena using polynomial approximations in an  $hp$  setting; we encounter row dependency in the experiments described in Section 3 and describe appropriate increases in  $[m]$ .

The set of acceptable arrays of coefficients  $\mathbf{b}$  that can be used to define functions in  $G_0[n][m]$  is the null space of  $\mathbf{M}$ , which we obtain as follows:

Compute the “QR” factorization of  $\mathbf{M}^T$ , so  $\mathbf{M}^T = (\mathbf{Q}'^T | \mathbf{Q}'^T) \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}$ , where  $\mathbf{R}$  is square upper-triangular

and invertible and  $\mathbf{Q} \equiv (\mathbf{Q}'^T | \mathbf{Q}'^T)$  is orthogonal. Then  $\mathbf{M} = \mathbf{R}^T \mathbf{Q}^T$ , where  $\mathbf{R} \equiv \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix}$ . Since we are looking for  $\mathbf{b}$  such that  $\mathbf{M}\mathbf{b}^T = \mathbf{0}$ , a straightforward argument shows that the columns of  $\mathbf{Q}'$ , the last  $m - n'$  columns of  $\mathbf{Q}$ , are an orthonormal basis for the null space of  $\mathbf{M}$ .

Let  $p \equiv m - n'$  and suppose that the  $p$  columns of  $\mathbf{Q}'$  are  $(q_{11}, \dots, q_{m1})^T, (q_{12}, \dots, q_{m2})^T, \dots, (q_{1p}, \dots, q_{mp})^T$ .

We enumerate the  $\{B_k^i\}$  as

$$\{B_1^1, B_2^1, \dots, B_{m_1}^1, B_1^2, B_2^2, \dots, B_{m_2}^2, \dots\};$$

there are  $m$  such  $B_k^i$ . Denote the  $B_k^i$  with this enumeration as  $\{\Phi^1, \Phi^2, \dots, \Phi^m\}$  and form  $\mathcal{B}_i \equiv \sum_{j=1}^m q_{ji} \Phi^j$  defined on all of  $\Omega$  by assuming each  $B_k^j$  is zero outside  $\Omega_j$ . Then  $\{\mathcal{B}_i\}$  is a basis for  $G_0[n][m]$ . Any approximation of form  $u_{n,m} = \sum_{i=1}^p \beta_i \mathcal{B}_i$  can be expressed in terms of the original basis represented by  $\{\Phi^1, \Phi^2, \dots, \Phi^m\}$  using coefficients  $\phi_j$ , the components of vector  $\phi = \mathbf{Q}'^T \mathbf{y}^T$ , where  $\mathbf{y} = (\beta_1, \beta_2, \dots, \beta_p)$ .

We consider the following non-self-adjoint problem: Let  $A$  be the linear operator

$$Au \equiv - \sum_{i,j=1}^K D_i(a_{ij} D_j u) + \sum_{i=1}^K a_i D_i u + a_0 u.$$

We consider the homogeneous Dirichlet problem  $Au = f$  with  $\gamma(u) = 0$ .

We treat elliptic problems, so we assume that the operator

$$Eu \equiv - \sum_{i,j=1}^K D_i(a_{ij}(x) D_j u) + a_0 u$$

is elliptic, self-adjoint and coercive. It suffices that the following hold:

(E1) All  $a_{ij}$  and  $a_0$  are in  $C^1(\overline{\Omega})$ ,

(E2) There exists  $c > 0$  such that

$$\sum_{i,j} a_{ij}(x) z_i z_j \geq c \sum_{i=1}^K z_i^2$$

for  $x$  in  $\Omega$  and for any  $z_i \in \mathbb{R}$  and  $a_{ij}(x) = a_{ji}(x)$  and  
 (E3)  $a_0(x) > 0$  for  $x \in \Omega$ .

We assume that the functions  $a_i$  forming the first-order term  $\sum_{i=1}^K a_i D_i u$ , denoted by  $\mathbf{a} \cdot \nabla u$ , are measurable and bounded.

Operator  $E$  induces a coercive bilinear form

$$a(u, v) \equiv \int_{\Omega} \sum_{i,j}^K a_{ij}(x) D_i u D_j v + a_0(x) uv \, dx.$$

There is a constant depending on the  $a_{ij}$  and  $a_0$  such that for any  $v \in H$ ,  $a(v, v) \leq M \|v\|_H^2$ .

Green's formula relates  $E$  and  $a(\cdot, \cdot)$ .

Let  $D_{\mathbf{n}_{ij}} u$  be the “co-normal derivative with respect to  $E$  of  $u$  on  $\Gamma_{ij}$ ” defined as follows: If  $\mathbf{n} = (n_1, n_2, \dots, n_K)$  is the unit normal to  $\Gamma_{ij}$  (pointing outward relative to the interior of  $\Omega_i$ ), then

$$D_{\mathbf{n}_{ij}} u \equiv \sum_{p,q}^K \gamma_{ij} (a_{pq} D_q u) n_p.$$

Let  $D_n u$  denote the general co-normal derivative. Green's formula

$$(\mathbf{E}u, v) = a(u, v) - \langle D_n u, \gamma(v) \rangle_{\Gamma'} \quad (3)$$

holds for suitable domains  $\Omega'$  with boundary  $\Gamma'$  for  $u$  in  $H^2(\Omega')$  and  $v$  in  $H^1(\Omega')$  if the  $a_{ij}$  are sufficiently smooth [19]. In particular, Green's formula is valid with our assumptions concerning the  $a_{ij}$  and  $\Omega' = \Omega$  or  $\Omega' = \text{any } \Omega_j$ .

The trace of the solution  $u$  is to be zero on  $\partial\Omega$ . The weak form of the problem is the following: For solution  $u \in H_0^1(\Omega)$ , we require that

$$a(u, v) + (\mathbf{a} \cdot \nabla u, v) = (f, v) \quad \text{for all } v \in H_0^1(\Omega).$$

To obtain an approximation, we use the Galerkin method with the basis  $\{\mathcal{B}_i\}$  for the space  $G_0[n][m]$ . An approximation is of form  $u_p \equiv \sum_{i=1}^p \beta_i \mathcal{B}_i(x)$ , and we require that

$$a(u_p, \mathcal{B}_j) + (\mathbf{a} \cdot \nabla u_p, \mathcal{B}_j) = (f, \mathcal{B}_j) \quad \text{for } j = 1, \dots, p.$$

We express this in terms of the original basis  $\{B_i^k\}$  as follows: Let  $\mathbf{C}$  denote the matrix of positive definite diagonal blocks  $(a(B_i^k, B_j^k))$ , so  $a(\mathcal{B}_i, \mathcal{B}_j) = \mathbf{Q}^T \mathbf{C} \mathbf{Q}$ . If we denote the (nonsymmetric) matrix of diagonal blocks  $((\mathbf{a} \cdot \nabla B_r^k, B_s^k))$  by  $\mathbf{H}$ ,  $((\mathbf{a} \cdot \nabla \mathcal{B}_i, \mathcal{B}_j)) = \mathbf{Q}^T \mathbf{H} \mathbf{Q}$ . Finally,  $(f, \mathcal{B}_j) = ((f, B_i^k)) \mathbf{Q}$ .

We have represented the row of undetermined coefficients by  $\beta$ ; we obtain our approximation by solving

$$\beta \mathbf{Q}^T (\mathbf{C} + \mathbf{H}) \mathbf{Q} = ((f, B_i^k)) \mathbf{Q}. \quad (4)$$

The size of the matrix  $\mathbf{Q}^T (\mathbf{C} + \mathbf{H}) \mathbf{Q}$  is  $p \times p$ , where

$p = (\text{size of matrix } \mathbf{C}) - (\text{total number of collocations enforced})$ .

For a problem with 16 rectangular cells, using the full basis for the set of polynomials of degree 10 or less and 8 collocations on each  $\Gamma_{ij}$ ,  $p$  is about 736. The  $p$  degrees of freedom in the system pertain to the approximation of the solution of the equation; we have eliminated any concern with weak continuity between interfaces of cells.

We give two convergence results.

If we assume that the problem is *elliptic* and *coercive* in the sense of Wloka [19], then there exist positive constants  $c_1$  and  $c_2$ , such that for any  $u$  and  $v$  in  $G_0[n]$ ,

$$(a) \quad |a(u, v) + (\mathbf{a} \cdot \nabla u, v)| \leq c_1 \|u\|_H \|v\|_H \text{ and}$$

$$(b) \quad a(u, u) + (\mathbf{a} \cdot \nabla u, u) \geq c_2 \|u\|_H^2.$$

The coercivity in (b) implies that the linear system has a unique solution  $\beta$ , which we use to define an approximate solution  $u_{n,m} \equiv \sum_{i=1}^p \beta_i \mathcal{B}_i(x)$ . Our first theorem makes the assumption that (a) and (b) hold. The second theorem assumes only that the system of linear equations (4) has a unique solution.

**Theorem 1.2.** Suppose that operator  $E$  satisfies (E1)–(E3), functions  $a_i$  are bounded and measurable,  $f \in L_2(\Omega)$ , and that  $u$  is a solution in  $H_0^2(\Omega)$  to  $Eu + \mathbf{a} \cdot \nabla u = f$ . Suppose that there is some  $c_2 > 0$ , such that for all  $v \in G_0[n]$ ,  $a(v, v) + (\mathbf{a} \cdot \nabla v, v) \geq c_2 \|v\|_H^2$ . Let  $M_1 \equiv \sqrt{K} \sup\{|a_i(x)|\}$ , where the supremum is taken over all  $x \in \Omega$  and  $1 \leq i \leq K$ . If  $u_{n,m}$  is the approximation obtained by solving the system of linear equations, then

$$c_2 \|u - u_{n,m}\|_H \leq (M + M_1) K([n]) \|\mathcal{Q}_{[m]} u\|_H + n_f \sqrt{N} \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\},$$

where  $K([n])$  is the parameter of Lemma 1.1 and  $n_f$  is the maximum number of faces of any of the  $N$  cells.

**Proof.** The Galerkin approximation satisfies the relation

$$a(u_{n,m}, v) + (\mathbf{a} \cdot \nabla u_{n,m}, v) = (f, v) \text{ for all } v \in G_0[n][m].$$

Let  $\mathcal{P}_m^n v$  denote the projection of  $v \in G_0[n]$  onto  $G_0[n][m]$ .

$$\begin{aligned} & a(u - u_{n,m}, u - u_{n,m}) + (\mathbf{a} \cdot \nabla(u - u_{n,m}), u - u_{n,m}) \\ &= a(u - u_{n,m}, u - \mathcal{P}_m^n u) + (\mathbf{a} \cdot \nabla(u - u_{n,m}), u - \mathcal{P}_m^n u) \\ & \quad + a(u - u_{n,m}, \mathcal{P}_m^n u - u_{n,m}) + (\mathbf{a} \cdot \nabla(u - u_{n,m}), \mathcal{P}_m^n u - u_{n,m}) \\ &= a(u - u_{n,m}, u - \mathcal{P}_m^n u) + (\mathbf{a} \cdot \nabla(u - u_{n,m}), u - \mathcal{P}_m^n u) \\ & \quad + a(u, \mathcal{P}_m^n u - u_{n,m}) + (\mathbf{a} \cdot \nabla u, \mathcal{P}_m^n u - u_{n,m}) \\ & \quad - a(u_{n,m}, \mathcal{P}_m^n u - u_{n,m}) - (\mathbf{a} \cdot \nabla u_{n,m}, \mathcal{P}_m^n u - u_{n,m}). \end{aligned}$$

Let  $\delta = \mathcal{P}_m^n u - u_{n,m} = \mathcal{P}_m^n(u - u_{n,m}) \in G_0[n][m]$ . Then

$$\begin{aligned} & a(u_{n,m}, \mathcal{P}_m^n u - u_{n,m}) + (\mathbf{a} \cdot \nabla u_{n,m}, \mathcal{P}_m^n u - u_{n,m}) \\ &= a(u_{n,m}, \delta) + (\mathbf{a} \cdot \nabla u_{n,m}, \delta) = (f, \delta). \end{aligned}$$

As in [17], we use Green's formula to get

$$\begin{aligned} & a(u, \mathcal{P}_m^n u - u_{n,m}) + (\mathbf{a} \cdot \nabla u, \mathcal{P}_m^n u - u_{n,m}) = a(u, \delta) + (\mathbf{a} \cdot \nabla u, \delta) \\ &= (Eu, \delta) + (\mathbf{a} \cdot \nabla u, \delta) + \sum_{I_{ij}} \langle D_{n_{ij}} u, \gamma_{ij}(\delta) - \gamma_{ji}(\delta) \rangle_{ij} + \sum_{I_{i0}} \langle D_{n_{i0}} u, \gamma_{i0}(\delta) \rangle_{i0} \\ &= (f, \delta) + \sum_{I_{ij}} \langle D_{n_{ij}} u, \gamma_{ij}(\delta) - \gamma_{ji}(\delta) \rangle_{ij} + \sum_{I_{i0}} \langle D_{n_{i0}} u, \gamma_{i0}(\delta) \rangle_{i0}. \end{aligned}$$

Using the argument in [17], since  $\delta$  is in  $G_0[n][m]$ , the differences of the traces of  $\delta$  on any  $\Gamma_{ij}$  are orthogonal to the  $[n]$  collocation functions and we can use Schwarz' inequality to obtain estimate

$$\left| \sum_{\Gamma_{ij}} \langle D_{n_{ij}} u, \gamma_{ij}(\delta) - \gamma_{ji}(\delta) \rangle_{ij} + \sum_{\Gamma_{i0}} \langle D_{n_{i0}} u, \gamma_{i0}(\delta) \rangle_{i0} \right| \\ \leq \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\} (n_f) \sqrt{N} \|\delta\|_H,$$

where  $C_{ij}$  are the trace constants,  $\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}$  is the norm of that part of the trace of the normal derivative of  $u$  on  $\Gamma_{ij}$  that is *not* in the span of the moment collocation functions,  $n_f$  is the maximum number of faces of any cell, and  $N$  is the number of cells. Recalling the definition of  $\delta$  we have

$$\|\delta\|_H = \|\mathcal{P}_m^n(u - u_{n,m})\|_H \leq \|u - u_{n,m}\|_H.$$

For any  $u, v \in H$ , the Schwarz inequality gives the estimate

$$|(a \cdot \nabla u, v)| \leq \sup\{|a_i(x)|\} \sum_{i=1}^K \|D_i u\|_0 \|v\|_0 \leq \sup\{|a_i(x)|\} k^{1/2} \|u\|_H \|v\|_0.$$

We have defined  $M_1$  to be  $\sup\{|a_i(x)|\} k^{1/2}$ . Recall that  $|a(u, v)| \leq M \|u\|_H \|v\|_H$ . These estimates give

$$\begin{aligned} & a(u - u_{n,m}, u - u_{n,m}) + (a \cdot \nabla(u - u_{n,m}), u - u_{n,m}) \\ & \leq M \|u - u_{n,m}\|_H \|u - \mathcal{P}_m^n u\|_H + M_1 \|u - u_{n,m}\|_H \|u - \mathcal{P}_m^n u\|_0 \\ & \quad + (f, \delta) + \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\} (n_f) \sqrt{N} \|u - u_{n,m}\|_H - (f, \delta) \\ & \leq \|u - u_{n,m}\|_H \{(M + M_1) \|u - \mathcal{P}_m^n u\|_H + \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\} n_f \sqrt{N}\}. \end{aligned} \quad (5)$$

By assumption

$$c_2 \|u - u_{n,m}\|_H^2 \leq a(u - u_{n,m}, u - u_{n,m}) + (a \cdot \nabla(u - u_{n,m}), u - u_{n,m}),$$

so term  $\|u - u_{n,m}\|_H$  can be canceled from both sides of the inequality above.

Since  $\|u - \mathcal{P}_m^n u\|_H \leq K([n]) \|\mathcal{Q}_{[m]} u\|_H$ , we obtain the estimate of the theorem.  $\square$

The coercivity condition  $a(u, u) + (a \cdot \nabla u, u) \geq c_2 \|u\|_H^2$  is stronger than necessary to insure existence of a unique solution [11, 19]. Grisvard [11], for example, shows that it suffices that  $a_0(x) \geq c_5 > 0$  and that the  $a_i$  be bounded and measurable. Our next theorem gives a result for approximations obtained under the assumption that there is a unique solution  $u_{n,m}$  to (4) and only form  $a(u, v)$  is coercive.

**Theorem 1.3.** Suppose that form  $a(u, v)$  satisfies the requirements of Theorem 1.2 and functions  $a_i$  are in  $C^1(\bar{\Omega})$  and there is some  $c_3 > 0$  such that

$$a_0(x) - \left(\frac{1}{2}\right) \nabla \cdot a(x) \geq c_3$$



for all  $x$  in  $\Omega$ . Then there is a constant  $c_4 > 0$  such that

$$\begin{aligned} c_4 \|u - u_{n,m}\|_H &\leq (M + M_1) K([n]) \|\mathcal{Q}_{[m]} u\|_H + \sqrt{N} n_f \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\} \\ &\quad + (\tfrac{1}{2}) \sqrt{n_f} \sup\{C_{ij}\} M_1 \left\{ \sum_{\Gamma_{ij}} \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij}^2 + \sum_{\Gamma_{i0}} \|\gamma_{i0}(u_{n,m})\|_{i0}^2 \right\}^{1/2}. \end{aligned}$$

**Proof.** For any  $v \in H$ ,

$$\begin{aligned} (\mathbf{a} \cdot \nabla v, v) &= \int_{\Omega} \left( \sum_{i=1}^K a_i D_i v \right) v \, dx = \frac{1}{2} \int_{\Omega} \sum_{i=1}^K a_i D_i v^2 \, dx \\ &= \frac{1}{2} \int_{\Omega} \nabla \cdot (v^2 \mathbf{a}) \, dx - \frac{1}{2} \int_{\Omega} v^2 \nabla \cdot \mathbf{a} \, dx. \end{aligned}$$

By the divergence theorem

$$\int_{\Omega} \nabla \cdot (v^2 \mathbf{a}) \, dx = \sum_{i=1}^N \int_{\Omega_i} \nabla \cdot (v^2 \mathbf{a}) \, dx = \sum_{i=1}^N \left( \sum_j \int_{\Gamma_{ij}} \gamma_{ij}(v^2 \mathbf{a}) \cdot \mathbf{n}_{ij} \, ds \right),$$

where  $\mathbf{n}_{ij}$  is the unit outward normal to  $\Gamma_{ij}$ .

Using the fact that if the boundary segment is an internal interface,  $\gamma_{ij}(\mathbf{a}) \cdot \mathbf{n}_{ij} = -\gamma_{ji}(\mathbf{a}) \cdot \mathbf{n}_{ji}$ , and grouping the sum of the boundary integrals above in pairs (if the boundary segment is an internal interface), we obtain

$$\begin{aligned} \sum_{i=1}^N \left( \sum_j \int_{\Gamma_{ij}} \gamma_{ij}(v^2 \mathbf{a}) \cdot \mathbf{n}_{ij} \, ds \right) &= \sum_{\Gamma_{ij}} \int_{\Gamma_{ij}} [\gamma_{ij}(v^2) - \gamma_{ji}(v^2)] \gamma_{ij}(\mathbf{a}) \cdot \mathbf{n}_{ij} \, ds \\ &\quad + \sum_{\Gamma_{i0}} \int_{\Gamma_{i0}} \gamma_{i0}(v^2) \gamma_{i0}(\mathbf{a}) \cdot \mathbf{n}_{i0} \, ds. \end{aligned}$$

The first sum in the expression above is taken over  $j \neq 0$  and we assume that  $i < j$ .

We let  $v = u - u_{n,m}$ . For internal interfaces  $\Gamma_{ij}$ ,

$$\begin{aligned} &\gamma_{ij}((u - u_{n,m})^2) - \gamma_{ji}((u - u_{n,m})^2) \\ &= \gamma_{ij}(u^2 - 2uu_{n,m} + u_{n,m}^2) - \gamma_{ji}(u^2 - 2uu_{n,m} + u_{n,m}^2) \\ &= \gamma_{ij}(u^2) - \gamma_{ji}(u^2) - 2\gamma_{ij}(u)[\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})] + [\gamma_{ij}(u_{n,m})]^2 - [\gamma_{ji}(u_{n,m})]^2. \end{aligned}$$

Since

$$\begin{aligned} \gamma_{ij}(u) &= \gamma_{ji}(u), \gamma_{ij}((u - u_{n,m})^2) - \gamma_{ji}((u - u_{n,m})^2) \\ &= [\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})][ -2\gamma_{ij}(u) + \gamma_{ij}(u_{n,m}) + \gamma_{ji}(u_{n,m}) ] \\ &= [\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})][\gamma_{ij}(u_{n,m} - u) + \gamma_{ji}(u_{n,m} - u)]. \end{aligned}$$

The homogeneous boundary condition gives  $\gamma_{i0}(u)=0$ , so

$$\gamma_{i0}((u - u_{n,m})^2) = \gamma_{i0}((u_{n,m})^2) = \gamma_{i0}(u_{n,m})\gamma_{i0}(u_{n,m} - u).$$

Assembling these results, we get

$$\begin{aligned} & 2(\mathbf{a} \cdot \nabla(u - u_{n,m}), u - u_{n,m}) \\ &= - \int_{\Omega} (u - u_{n,m})^2 \nabla \cdot \mathbf{a} \, dx + \sum_{\Gamma_{ij}} \int_{\Gamma_{ij}} [\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})][\gamma_{ij}(u_{n,m} - u) \\ & \quad + \gamma_{ji}(u_{n,m} - u)]\gamma_{ij}(\mathbf{a}) \cdot \mathbf{n}_{ij} \, ds + \sum_{\Gamma_{i0}} \int_{\Gamma_{i0}} \gamma_{i0}(u_{n,m})\gamma_{i0}(u_{n,m} - u)\gamma_{i0}(\mathbf{a}) \cdot \mathbf{n}_{i0} \, ds. \end{aligned}$$

The use of Schwarz' inequality readily shows that constant  $M_1$  of Theorem 1.2 majorizes the supremum of  $|\gamma_{ij}(\mathbf{a}) \cdot \mathbf{n}_{ij}|$  over all  $\Gamma_{ij}$  and  $\Gamma_{i0}$ . Then

$$\begin{aligned} & \left| \int_{\Gamma_{ij}} [\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})][\gamma_{ij}(u_{n,m} - u) + \gamma_{ji}(u_{n,m} - u)]\gamma_{ij}(\mathbf{a}) \cdot \mathbf{n}_{ij} \, ds \right| \\ & \leq M_1 \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij} [\|\gamma_{ij}(u_{n,m} - u)\|_{ij} + \|\gamma_{ji}(u_{n,m} - u)\|_{ji}] \\ & \leq M_1 \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij} [C_{ij}\|u_{n,m} - u\|_{1,i} + C_{ji}\|u_{n,m} - u\|_{1,j}]. \end{aligned}$$

Likewise

$$\left| \int_{\Gamma_{i0}} \gamma_{i0}(u_{n,m})\gamma_{i0}(u_{n,m} - u)\gamma_{i0}(\mathbf{a}) \cdot \mathbf{n}_{i0} \, ds \right| \leq M_1 \|\gamma_{i0}(u_{n,m})\|_{i0} C_{i0} \|u_{n,m} - u\|_{1,i}.$$

Using the Schwarz inequality, and the fact that any  $\|u_{n,m} - u\|_{1,i}$  occurs at most  $n_f$  times in the sums, the sums over the  $\Gamma_{ij}$  and  $\Gamma_{i0}$  are majorized by

$$M_1 \left\{ \sum_{\Gamma_{ij}} \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij}^2 + \sum_{\Gamma_{i0}} \|\gamma_{i0}(u_{n,m})\|_{i0}^2 \right\}^{1/2} \sup\{C_{ij}\} (n_f)^{1/2} \left\{ \sum_{i=1}^N \|u_{n,m} - u\|_{1,i}^2 \right\}^{1/2}.$$

This last sum is  $\|u_{n,m} - u\|_H$ . We combine these estimates with inequality (5) of Theorem 1.2:

$$\begin{aligned} & a(u - u_{n,m}, u - u_{n,m}) - \frac{1}{2} \int_{\Omega} (u - u_{n,m})^2 \nabla \cdot \mathbf{a} \, dx \\ & \leq \|u - u_{n,m}\|_H \left[ (M + M_1) \|u - \mathcal{P}_m^n u\|_H + \sup\{C_{ij}\} \sup\{\|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}} u)\|_{ij}\} n_f \sqrt{N} \right. \\ & \quad \left. + \left( \frac{1}{2} \right) \sup\{C_{ij}\} (n_f)^{1/2} M_1 \left\{ \sum_{\Gamma_{ij}} \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij}^2 + \sum_{\Gamma_{i0}} \|\gamma_{i0}(u_{n,m})\|_{i0}^2 \right\}^{1/2} \right] \end{aligned} \quad (6)$$

The supposition concerning  $c_3$  means that

$$\begin{aligned} a(u - u_{n,m}, u - u_{n,m}) - \frac{1}{2} \int_{\Omega} (u - u_{n,m})^2 \nabla \cdot \mathbf{a} \, dx \\ = \int_{\Omega} \sum_{i,j}^K a_{ij}(x) D_i(u - u_{n,m}) D_j(u - u_{n,m}) + \left[ a_0(x) - \left( \frac{1}{2} \right) \nabla \cdot \mathbf{a} \right] (u - u_{n,m})^2 \, dx \\ \geq c_4 \|u - u_{n,m}\|_H^2, \end{aligned}$$

so the term  $\|u - u_{n,m}\|_H$  can be canceled from each side of inequality (6) above.  $\square$

Due to the moment collocation constraints, the sums in the estimate over the  $\Gamma_{ij}$  and  $\Gamma_{i0}$  are small; more so as  $[n]$  is large. Such interface error terms are independent of the solution  $u$ , and are readily computed in the solution process. In the polynomial implementation in Section 3 approximations can be made continuous, thus eliminating these error terms.

If the vector field  $\mathbf{a}$  is constant or  $\nabla \cdot \mathbf{a} \leq 0$ , the assumptions concerning the existence of an appropriate constant  $c_3$  are unnecessary. In fact, a stronger result holds in this case; the following corollary also implies that the estimates of Theorem 1.3 apply to the steady-state convection–diffusion equations

$$\nabla \cdot (-\sigma \nabla u + u \mathbf{a}) = 0$$

for concentration  $u(x)$ , when the underlying flow with velocity  $\mathbf{a}$  is incompressible.

**Corollary 1.4.** *Suppose that form  $a(u, v)$  satisfies the requirements of Theorem 1.2 and functions  $a_i$  are in  $C^1(\bar{\Omega})$ . If the vector field  $\mathbf{a}$  satisfies  $\nabla \cdot \mathbf{a} \leq 0$ , and  $[n]$  is sufficiently large so that for each  $\Gamma_{ij}$  there is some  $k \leq n_{ij}$  such that  $\langle \omega_k^{ij}, 1 \rangle_{ij} \neq 0$ , then coefficient  $a_0(x)$  can be zero and the estimate of Theorem 1.3 is valid.*

**Proof.** From [4], when the additional assumption is made concerning the size of  $[n]$ , functional  $a(u, v)$  is coercive over  $G_0[n]$ , and the term containing nonpositive  $\nabla \cdot \mathbf{a}$  in Eq. (6) can be discarded.  $\square$

### 3. Nonhomogeneous boundary value problems

We adapt the algorithm to nonhomogeneous boundary value problems, where, given some  $g(\cdot)$  defined on  $\Gamma$ , we now require that  $u(x) = g(x)$ . The classical method is to first find some  $u_2(x)$  such that  $u_2 = g$  on  $\Gamma$ , next express  $u$  as  $u_1 + u_2$ , with  $u_1 = 0$  on  $\Gamma$ , and then solve  $Au_1 = f - Au_2$ .

The equivalent transformed problem requires that we find some vector  $\eta$  such that  $M\eta = g$ , where the components of  $g$  corresponding to internal collocation on  $\Gamma_{ij}$  are zero; collocation rows corresponding to  $\Gamma_{i0}$  produce an entry of form  $\langle g(\cdot), \omega_q^{i0} \rangle_{i0}$ . An  $\eta$  of minimal norm satisfying these requirements can be obtained as part of the **QR** factorization of  $M^T$ : if  $M^T = (Q'' | Q') \begin{pmatrix} R \\ 0 \end{pmatrix}$ , where

$\mathbf{R}$  is square upper-triangular and invertible, then  $\eta = \mathbf{Q}''(\mathbf{R}^T)^{-1}\mathbf{g}$ . We approximate  $u_2$  by

$$u_2|_{\Omega_k} = \sum_{i=1}^{m_k} \eta_i^k B_i^k,$$

where  $\eta_i^k$  is the component of  $\eta$  associated with  $B_i^k$ . Suppose that our approximation to  $u_1$  expressed in terms of the original basis is

$$u_1|_{\Omega_k} = \sum_{i=1}^{m_k} b_i^k B_i^k.$$

If we take the  $L_2(\Omega_k)$  inner product of  $Au_1 = f - Au_2$  with  $B_i^k$  and follow the argument in Section 2, the resulting vector equation is

$$\mathbf{b}(\mathbf{C} + \mathbf{H}) = \mathbf{f} - \eta^T(\mathbf{C} + \mathbf{H}),$$

where  $\mathbf{f} \equiv ((f, B_i^k)_0)$  and  $\mathbf{b}$  is to satisfy  $\mathbf{M}\mathbf{b}^T = \mathbf{0}$ .

In terms of the basis  $\{\mathcal{B}_j\}$  satisfying the collocation constraints,  $u_1$  is represented as  $\sum_i^p \beta_i \mathcal{B}_i(x)$  and vector  $\beta$  is to satisfy

$$\beta \mathbf{Q}'^T(\mathbf{C} + \mathbf{H})\mathbf{Q}' = \mathbf{f}\mathbf{Q}' - \eta^T(\mathbf{C} + \mathbf{H})\mathbf{Q}'.$$

We have implemented this scheme for arbitrary problems with domains in  $\mathbb{R}^2$  that can be partitioned into triangles or parallelograms (or both). We use  $L_2$ -orthonormal bases of polynomials of degree 10 or less (up to 66 functions on each cell) to provide approximations. Gauss–Legendre quadrature is used to effect the integrations over the cells and interfaces, and subroutines from LINPACK [5] and LAPACK [1] provide the  $QR$  decomposition and the solution of the final system.

We test the error estimates in terms of  $p$ , the degree of the polynomial approximation on each cell,  $q$ , the maximum degree of the Legendre polynomials providing the weight functions on the interfaces in any trial, and  $h$ , the maximum diameter of the cells in the cell decomposition of  $\Omega$ . We use the same number  $q + 1$  of collocations on each boundary segment  $\Gamma_{ij}$ , so we replace all  $n_{ij}$  with  $q + 1$ ; we revise the notation containing collocation index  $[n]$  by replacing  $[n]$  with  $q$ . We also use the same number of basis functions for our  $p$ th order basis on each cell; the notation containing basis multi-index  $[m]$  replaces  $[m]$  by  $p$ . Thus, we now denote approximation  $u_{n,m}$  with symbol  $u_{q,p}$ . The relevant error estimates for a polynomial implementation of these methods are given in [4]. They are expressed in terms of  $H^k(\Omega)$  norm of the solution ( $k > 2$ ) or, for analytic solutions, in terms of the semi-norm defined by the  $L_2$ -norm of the  $p + 1$  and  $p + 2$  derivatives of the solution. Since our test problem is analytical we use these second estimates.

The following estimates hold for domains in  $\mathbb{R}^2$ :

The trace constants  $C_{ij}$  for boundaries of polygonal cells are bounded by  $c_1/(h)^{1/2}$ , where  $c_1$  is independent of  $h$  and depends only on the smallest angle in any cell.

$K[n]$  is bounded by  $c_2(q)/h$ . We can get some sense of the size of  $c_2(q)$  from the experiments in [4]. For example, given even  $q$  and square cells,  $p \equiv q + 2$  is often optimal. In this case, we have computed  $c_2(q)$  explicitly for  $q = 2$ –20; a tight estimate is  $c_2(q) \cong 15q^{1.4}$ .

We assume that the coefficients  $a_{ij}$  and  $a_0$  are constant. Then

$$\|\mathcal{T}_q^{ij}(D_{n_{ij}}u)\|_{ij} \equiv \|\mathcal{T}_{n_{ij}}^{ij}(D_{n_{ij}}u)\|_{ij} \leq 0.66 \times h^{q+1} (0.7(q+2))^{-(q+3/2)} \|(D_{n_{ij}}u)^{q+1}\|_{ij},$$

where  $(D_{n_{ij}}u)^{q+1}$  represents the  $(q+1)$ st tangential derivative of  $D_{n_{ij}}u$  on  $\Gamma_{ij}$ .

For our cells, we have for  $v \in H^{p+2}(\Omega)$  (and for  $h \leq 3$  and  $p \geq 2$ ),

$$\|\mathcal{Q}_p(v)\|_H \leq h^p (0.5p)^{-p} [|v|_{p+1} + |v|_{p+2}], \quad (7)$$

where, for example,

$$|v|_{p+1}^2 = \sum_{|\alpha|=p+1} \|D^\alpha v\|_0^2.$$

If we use these estimates in Theorem 1.3., with  $q$  and  $p$  replacing  $[n]$  and  $[m]$ , and  $n_f \leq 4$ , we get

$$\begin{aligned} c_4 \|u - u_{q,p}\|_H &\leq (M + M_1) c_2(q) h^{p-1} (0.5p)^{-p} [|u|_{p+1} + |u|_{p+2}] \\ &\quad + 2.64 c_1 \sqrt{N} h^{q+1/2} (0.7(q+2))^{-(q+3/2)} \max \| (D_{n_{ij}}u)^{q+1} \|_{ij} \\ &\quad + c_1 h^{-1/2} M_1 \left\{ \sum_{\Gamma_{ij}} \|\gamma_{ij}(u_{q,p}) - \gamma_{ji}(u_{q,p})\|_{ij}^2 + \sum_{\Gamma_{i0}} \|\gamma_{i0}(u_{q,p})\|_{i0}^2 \right\}^{1/2}. \end{aligned} \quad (8)$$

If we are subdividing the unit square into cells of side  $h$ , the number of cells  $N \cong 1/h^2$ . Due to the resulting decrease in the size of  $\Gamma_{ij}$ , we might expect  $\|(D_{n_{ij}}u)^{q+1}\|_{ij}^2$  to decrease by a factor  $h$ . Then the  $h$ -dependency of the second error estimate term containing the normal derivative of the solution on the interfaces would be  $h^{-1} h^{q+1/2} h^{1/2} = h^q$ . This estimate can be made rigorous by the methods used in [4].

We report on two numerical tests of the theory.

**Example 1.** We approximate the solution to the following convection–diffusion problem from [14]:

$$-\sigma \Delta u + \mathbf{a} \cdot \nabla u = 0,$$

where  $\mathbf{a} = (a_1, a_2)$  is a constant unit vector giving the direction of the flow. The domain is an isosceles trapezoid with base on interval  $[0, 1]$  on the  $x$ -axis, acute angles of size  $\pi/4$  and top on the line  $y = 0.25$ . Vector  $\mathbf{a}$  points upward and makes an angle of  $75^\circ$  with the  $x$ -axis, so that the flow enters the domain across the  $x$ -axis. There is considerable interest in convection-dominated flows [12, 14]; we set  $\sigma$  to be 0.0125 in this first example. The boundary data is defined by the intended solution  $u(x, y) \equiv v(s(x, y), t(x, y))$ , where  $s = a_1 x + a_2 y$ ,  $t = a_2 x - a_1 y$ , and

$$v(s, t) = e^{-\lambda s} \sin(\omega t)$$

with  $\omega = 2\pi/a_2$  and  $\lambda = (\sqrt{1 + 4\omega^2 \sigma^2} - 1)/(2\sigma)$ .

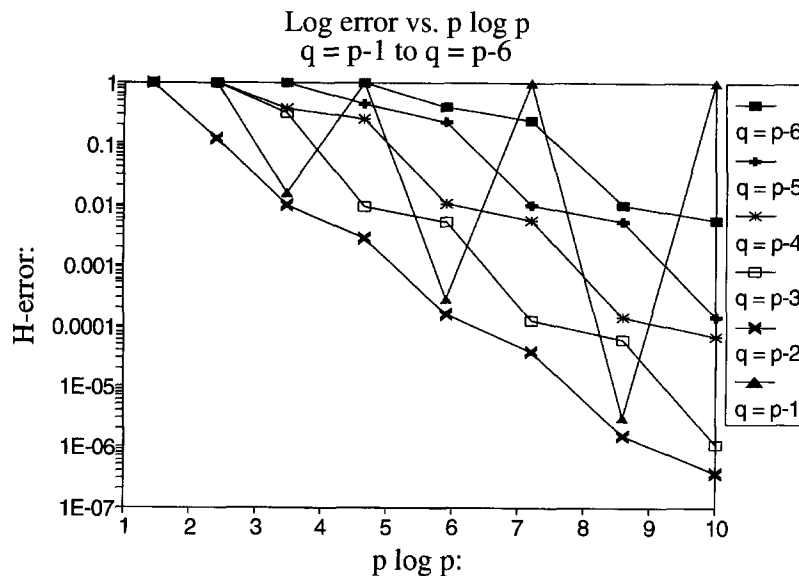


Fig. 1. Log of the  $H$ -norm of errors vs.  $p \log p$  for  $q = p - i$ .

Our first estimate is made on a single triangular cell containing the trapezoid; we subsequently partition the trapezoid into 4 and then 12 triangles, where the lengths of the base(s) of the triangle(s) are 1, 0.5 and 0.25. We study the effect of varying  $h, q$  and  $p$  on the accuracy of the approximation.

Fig. 1 plots the  $H$ -norm of the errors for  $p$  and  $q = p - i$ ,  $p = 3, \dots, 10$  and  $i$  ranging from 1 to 6 when the domain is partitioned into three triangles.

The approximation fails for even  $p$  when  $q = p - 1$  since here we have used more than enough collocations to force continuity; there is a dependency among the rows of  $M$ . The optimal value for  $q$  is  $p - 2$  (when we use a partition into triangles); thus, requiring that the approximation be continuous does not seem to give any advantage. In all cases we computed the interface error

$$\left\{ \sum_{I_{ij}} \|\gamma_{ij}(u_{n,m}) - \gamma_{ji}(u_{n,m})\|_{ij}^2 + \sum_{I_{i0}} \|\gamma_{i0}(u_{n,m})\|_{i0}^2 \right\}^{1/2},$$

which is part of the general error estimate; in our tests the interface error is less than  $\frac{1}{20}$  of the computed error. For example, when  $p = 10$  and  $q = 8$ , the computed error is  $0.36 \times 10^{-6}$  while the interface error is only  $0.64 \times 10^{-8}$ . The regression equation for the error when  $q = p - 2$  is

$$H\text{-error} \cong 0.94(0.46p)^{-p}.$$

This agrees well with the estimate in (8).

The  $h$ -dependency for the error shown is Fig. 2, where  $q$  is set to  $p - 2$ . The approximate slopes of the lines are shown on the right; they are about  $q + \frac{1}{2}$ ; they are the powers of  $h$  in our regression

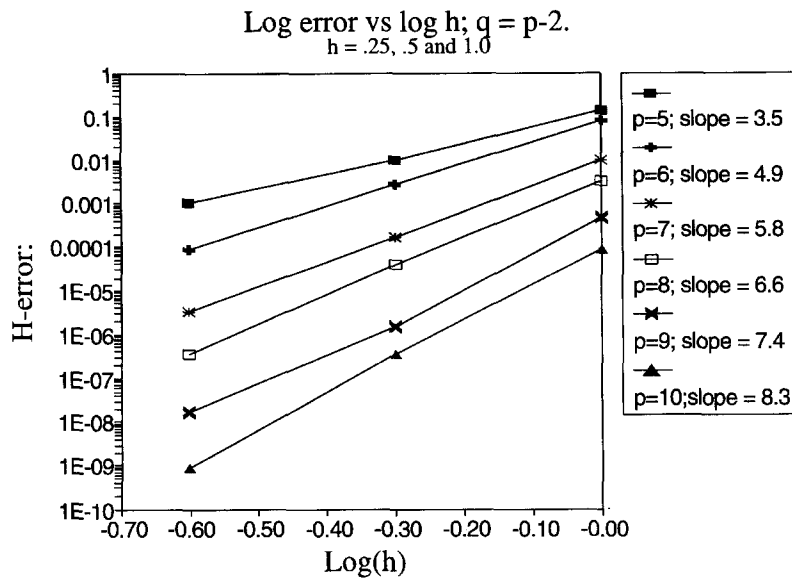


Fig. 2. Log error vs. log  $h$ ;  $q = p - 2$ ;  $h = 0.25, 0.5$  and  $1$ .

model for the error of form  $H\text{-error} \cong Ch^s$ ; thus, our experiments give error of form  $Ch^{q+0.5}$ , which is consistent with our theory (8).

**Example 2.** Traditional Galerkin methods can encounter difficulty when there is a boundary layer that can form with such convection-dominated flows [13]. We study such a problem here. Our equation is

$$-\sigma \Delta u + u_y = - \left( \frac{1}{2} + 2\sigma \pi^2 (y + 1 + C_1) \right) \sin(2\pi x)$$

with domain  $(0, 1) \times (0, 1)$  and  $C_1 \equiv 2[\exp(1/\sigma) \sinh((1 + 16\sigma^2 \pi^2)^{1/2}/\sigma) - 1]^{-1} \cong 4 \exp(-2/\sigma)$ , negligible for small  $\sigma$ . Dirichlet boundary data is zero on  $y = 1$  and  $x = 0$  and  $x = 1$ ; on  $y = 0$ , the inflow portion of  $\Gamma$ , the boundary values are essentially  $-0.5 \sin(2\pi x)$  (with small sigma); more precisely, they are given by the intended solution

$$u(x, y) = \left[ \frac{\exp((y + 1)/2\sigma) \sinh(\sqrt{1 + 16\sigma^2 \pi^2} (y + 1)/2\sigma) - 1}{\exp(1/\sigma) \sinh(\sqrt{1 + 16\sigma^2 \pi^2}/\sigma) - 1} - \frac{y + 1}{2} \right] \sin(2\pi x),$$

evaluated at  $y = 0$ . There is a boundary layer at  $y = 1$  of width  $O(\sigma \ln(2/\sigma))$ .

Fig. 3 shows two approximations when  $\sigma = 0.005$ , with  $p = 10$  and  $q = 7$ . In Fig. 3(a) we partition the unit square into 4 congruent rectangles with base 0.25 and height 1. The boundary layer causes considerable distortion in the approximation, causing the un-physical oscillation typical of such problems [12, 14]. The width of the boundary layer is about  $0.03 \cong 0.005 \ln(2/0.005)$ . In Fig. 3(b) we show the approximation when we use 8 cells, where each of the 4 cells above is partitioned into two

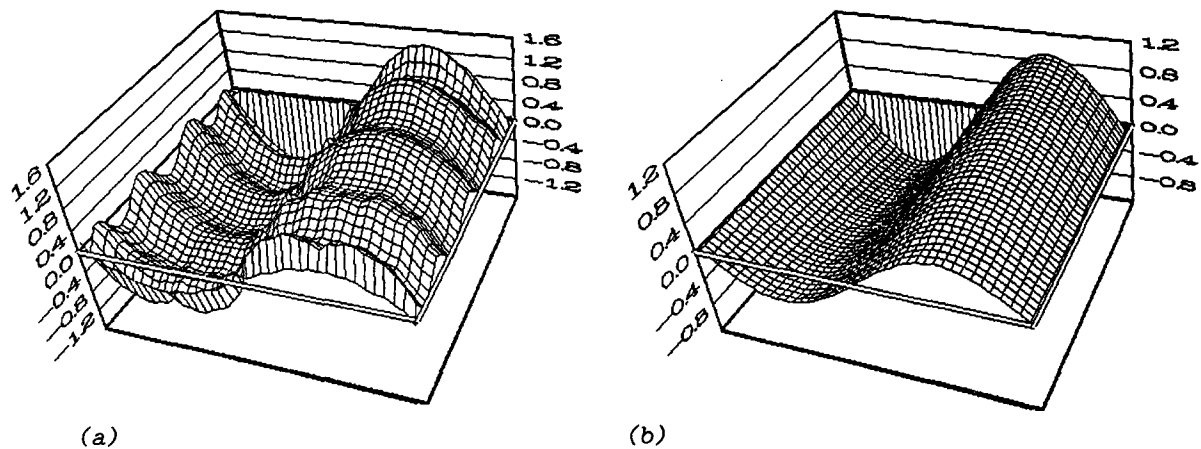


Fig. 3. Approximations to the boundary layer problem.

rectangles; a large one of width 0.25 and height 0.97 and a thin one of height 0.03 to accommodate the boundary layer. The result in Fig. 3(b) has maximum error less than 0.001.

The choice of optimal meshes for a similar problem with boundary layers when using the  $hp$  finite element method is treated in [16]; these results suggest directions for further investigations of the methods we have discussed in this paper.

## References

- [1] E. Anderson et al., LAPACK Users' Guide, SIAM, Philadelphia, PA, 1991.
- [2] I. Babuška, M.R. Dorr, Error estimates for the combined  $h$  and  $p$  versions of the finite element method, *Numer. Math.* 37 (2) (1981) 257–277.
- [3] I. Babuška, M. Suri, The  $p$  and  $h$ - $p$  versions of the finite element method, an overview, *Comput. Meth. Appl. Mech. Eng.* 80 (1–3) (1990) 5–26.
- [4] M. Cayco, L. Foster, H. Swann, On the convergence rate of the cell discretization algorithm for solving elliptic problems, *Math. Comp.* 64 (212) (1995) 1397–1419.
- [5] J. Dongarra et al., LINPACK Users Guide, SIAM, Philadelphia, PA, 1979.
- [6] M.R. Dorr, On the discretization of interdomain coupling in elliptic boundary-value problems, in: T.F. Chan, R. Glowinski, J. Periaux, O.B. Widlund (Eds.), *Domain Decomposition Methods*, SIAM, Philadelphia, PA, 1989.
- [7] J. Greenstadt, Cell discretization, in: J.H. Morris (Ed.), *Conference on Applications of Numerical Analysis, Lecture Notes in Mathematics*, 228, Springer, New York, 1971, pp. 70–82.
- [8] J. Greenstadt, The cell discretization algorithm for elliptic partial differential equations, *SIAM J. Sci. Statist. Comput.* 3 (3) (1982) 261–288.
- [9] J. Greenstadt, The application of cell discretization to nuclear reactor problems, *Nucl. Sci. Eng.* 82 (1982) 78–95.
- [10] J. Greenstadt, Cell discretization of nonselfadjoint linear elliptic partial differential equations, *SIAM J. Sci. Stat. Comput.* 12 (5) (1991) 1074–1108.
- [11] P. Grisvard, *Elliptic problems in non-smooth domains*, Pitman, UK, 1985.
- [12] B. Guo, I. Babuška, The  $h$ - $p$  version of the finite element method. Part 1: The basic approximation results. Part 2: General results and applications, *Comput. Mech.* 1 21–41 (1986) 203–226.
- [13] W.B. Liu, J. Shen, A new efficient spectral Galerkin method for singular perturbation problems, *J. Sci. Comput.* 11 (1996) 411–437.



- [14] P.A. Raviart, J.M. Thomas, Primal hybrid finite element methods for second order elliptic equations, *Math. Comp.* 31 (138) (1977) 391–413.
- [15] G.R. Richter, An explicit finite element method for convection-dominated steady state convection–diffusion equations, *SIAM J. Numer. Anal.* 28 (3) (1991) 744–759.
- [16] C. Schwab, M. Suri, The  $p$  and  $hp$  versions of the finite element method for problems with boundary layers, *Math. Comp.* 1997, to appear.
- [17] H. Swann, On the use of Lagrange multipliers in domain decomposition for solving elliptic problems, *Math. Comp.* 60 (201) (1993) 49–78.
- [18] H. Swann, Error estimates using the cell discretization method for some parabolic problems, *J. Comput. Appl. Math.* 66 (1996) 479–514.
- [19] J. Wloka, *Partial Differential Equations*, Cambridge University Press, Cambridge, 1987.